

HEPScore: A new CPU benchmark for the WLCG

Domenico Giordano¹, Jean-Michel Barbe², Tommaso Boccali³, Gonzalo Menéndez Borge⁴, Christopher Hollowell⁴, Vincenzino Innocenti⁵, Walter Lampf⁶, Michele Michelotto⁶, Helge Meinhard⁷, Ladislav Ondris⁸, Andrea Sciabà⁸, Matthias J. Schnepf⁹, Randall J. Sobie⁹, David Southwick⁹, Tristan S. Sullivan⁹, Andrea Valassi¹⁰, Sandro Wenzel¹⁰, John L. Willis⁹, and Xiaofei Yan¹⁰

¹CERN, Geneva, Switzerland

²Subatech UMR 6457, CNRS-IN2P3, IMT Atlantique, Université de Nantes, Nantes, France

³INFN Sezione di Pisa, L.go Pontecorvo 3, 56127 Pisa (ITALY), and ICSC-National Center for HPC, Big Data and Quantum Computing, Italy

⁴Brookhaven National Laboratory, United States

⁵University of Arizona, United States

⁶INFN, Istituto Nazionale di Fisica Nucleare, Padua, Italy

⁷Karlsruhe Institute of Technology, Karlsruhe, Germany

⁸Institute of Particle Physics and University of Victoria, Victoria, British Columbia, Canada

⁹California Institute of Technology, Pasadena, California, United States

¹⁰Institute of High Energy Physics, Beijing, China

Abstract. HEPscore is a new CPU benchmark created to replace the HEP-SPEC06 benchmark that is currently used by the WLCG for procurement, computing resource pledges, usage accounting and performance studies. The development of the new benchmark, based on HEP applications or workloads, has involved many contributions from software developers, data analysts, experts of the experiments, representatives of several WLCG computing centres and WLCG site managers. In this contribution, we review the selection of workloads and the validation of the new HEPscore benchmark.

1 Introduction

Computing in particle physics has evolved to a highly distributed model where each country provides local facilities that are integrated into a global infrastructure, called the Worldwide LHC Computing Grid (WLCG) [1]. The WLCG coordinates the computing and networking resources on behalf of the experiments at the Large Hadron Collider at CERN.

The collaborative nature of our field has resulted in agreements for the sharing of costs for all aspects of our experiments, including computing, either through direct financial payments or the provision of in-kind equipment or services. Each country is requested to contribute resources based on the size of their research community and financial oversight boards find consensus on the appropriate cost sharing.

It is difficult to put a cost estimate on a computing facility as each country has its own way of acquiring and operating their resources. Rather than develop a model based on cost

e-mail: Domenico.Giordano@cern.ch

of the facility, hardware and personnel, it was decided to use the CPU-power delivered by a site to compare the CPU resources provided by each country. The delivered CPU-power is defined to be the number of seconds used by the applications multiplied by a benchmark that reflects the performance of the servers. After many years of operations, most sites have many types of servers with differing levels of performance. As a result, the sites report a benchmark that is average of the individual benchmarks of the different servers weighted by number of servers of each type (we refer to this number as *site benchmark*).

The resources used at each site are tracked and stored in a WLCG accounting database. The database stores the site benchmark and the number of CPU-seconds used by each experiment at a site. These numbers are used to calculate an integrated number that estimates the resources delivered by the site. These numbers are published by the WLCG accounting team on a monthly basis and provided to the funding agencies.

2 Motivation for a new CPU benchmark

In 2009, the WLCG agreed to use HEP-SPEC06 (HS06) as its CPU benchmark [2]. HS06 is based on the industry-standard SPEC CPU 2006 benchmark [3]. At that time, the HEP applications shared several commonalities with a number of workloads within the SPEC CPU 2006 suite and those workloads were selected to be included in HS06 [4, 5]. The individual HS06 workloads were single-threaded and single-process applications, compiled in 32-bit mode, and requiring a minimum of 1 GB of memory per process.

The support of the SPEC CPU 2006 benchmark has ended and the Working Group considered a number of options for replacement for HS06. A newer release of the SPEC benchmark, SPEC CPU 2017 [6], was also studied but found to be highly correlated with the HEP-SPEC06 benchmark (see fig. 1). The SPEC CPU 2017 benchmark would require each site to purchase a new licence (like the SPEC CPU 2006 benchmark) whereas the WLCG community was in favour an open-source benchmark. As a result, the SPEC CPU 2017 was rejected as a replacement for HS06.

HS06 has met the WLCG requirements in a world that has progressively evolved from CPUs with few cores to multi-core CPUs. During this period, experiments were starting to observe that newer versions of their applications scaled less well with HS06 (for example, see ref. [7]). The proposal for building a new benchmark using HEP applications was first presented at the 2017 WLCG workshop [8]. The HEPiX Benchmarking Working Group was asked to study the feasibility of a new benchmark based on HEP workloads. This led to the creation of the HEP Benchmarks project and the development of a software framework to benchmark the performance of computing servers using HEP applications [9]. The Working Group developed the HEP Benchmark Suite that provides a simple way to run containerized HEP applications and other benchmarks, and record the results in an Elasticsearch database at CERN [10].

Some of the HEP applications that use the largest amount of compute power are the Monte Carlo generation of collision events, the simulation of the detector response to the simulated particles, the conversion of the simulated energy deposition by the particles in the detector elements (digitization) and reconstruction of the detector signals into particles and momenta (the reconstruction code is used for both simulated and real data). We refer to these HEP applications as *workloads*. Analysis applications are very specific to the physics, and often I/O intensive, and are considered too unreliable to use as a measure of the performance of a CPU.

¹HEPiX Benchmarking Working Group: D Giordano (co-chair), M Michelotto (co-chair), L Atzori, JM Barbet, GM Borge, C Hollowell, L Ondris, A Sciaba, E Simili, RJ Sobie, D Southwick, T Sullivan, N Szczepanek, A Valassi

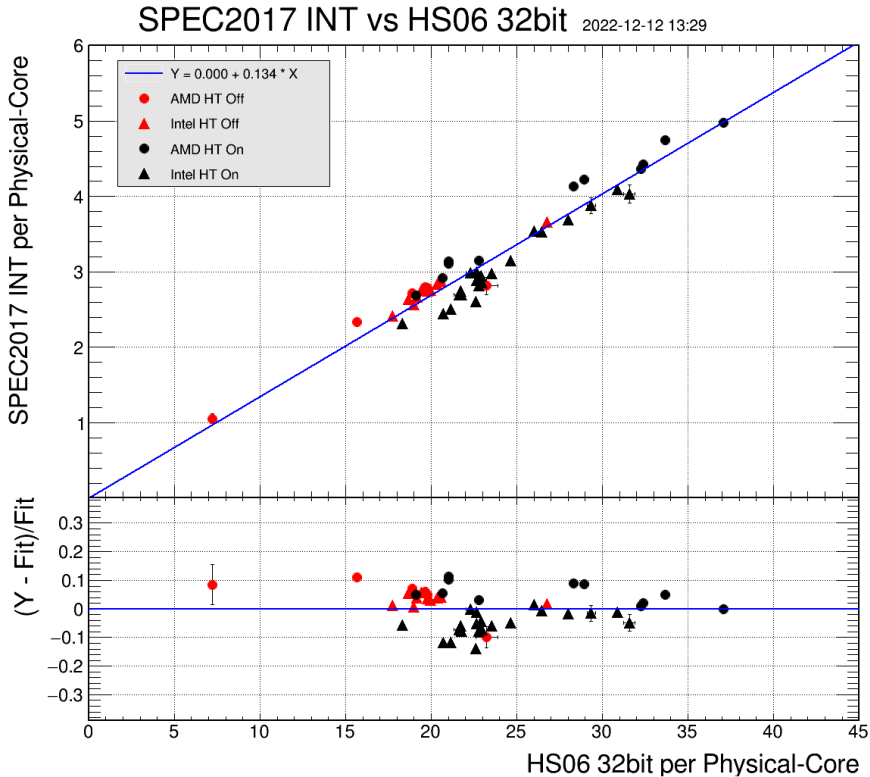


Figure 1. A comparison of the HEP-SPEC06 benchmark with the results of the SPEC CPU 2017 benchmark on resources provided to the HEPiX Benchmark Working Group. The circles are measurements on AMD processors and triangles are on Intel processors. The red (black) points have hyper-threading on (off). The blue line is a linear fit to the data constrained to the origin (0,0). The lower plot shows the fractional difference in the vertical dimension of the data point from the fit. Note that the benchmarks are normalized to the number of physical cores of the server independent of whether hyper-threading is enabled.

The creation of a new benchmark based on HEP workloads (called HEPscore) benefits from the consensus of the WLCG community. As a result, the WLCG Management Board established the HEPscore Benchmark Deployment Task Force whose role was to review the requirements for the HEPscore benchmark and to help select the HEP workloads that are to be used in the HEPscore benchmark. The Task Force was also asked to propose a transition plan for the migration from the HS06 benchmark to the new HEPscore benchmark.

In September 2022, a 2-day workshop at CERN devoted to the HEPscore benchmark brought together members of the Working Group, Task Force, computing coordinators and software developers of the experiments and site representatives to discuss the composition of the new benchmark [11]. A presentation at the 2022 ACAT conference provided an interim

²WLCG HEPscore Deployment Task Force: D Giordano (co-chair), RJ Sobie (co-chair), J Andreeva, G Andronico, GM Arevalo, T Boccali, GM Borge, C Bozzi, S Campana, I Collier, A Di Girolamo, M Jouvin, W Lampl, JR Letts, Z Marshall, H Meinhard, AM Melo, B Panzer-Steindel, S Piano, D Piparo, F Qi, MJ Schnepf, O Smirnova, J Templon, A Valassi, JL Willis, T Wong

report on the status of the HEPsScore benchmark [12]. The proposal for the first version of HEPsScore was presented at the 2022 WLCG Workshop [13] where a number of recommendations were proposed (discussed in the next section). A meeting of the WLCG Management Board reviewed the status of the new benchmark in December 2022 and set a milestone of April 2023 for the release of HEPsScore.

Concurrent to the development of the HEPsScore benchmark, the issue of power consumption and environmental impact of computing resources has become a topical area of interest [14]. A study, using a beta version of HEPsScore, evaluated the processing capabilities and power consumption of x86 and ARM processors, and showed that ARM processors use less power while still being highly performant for HEP workloads [15]. As a result, there was strong community interest in a benchmark for both x86 and ARM processors, and a concerted effort by the experiments resulted in all workloads being ready for both processor types for the April 2023 milestone.

3 HEP Workloads

In 2021 and 2022, workloads were provided to the Working Group by all four large LHC experiments (ALICE, ATLAS, CMS and LHCb) and other experiments (Belle II, JUNO and the International Gravitational Wave Network).

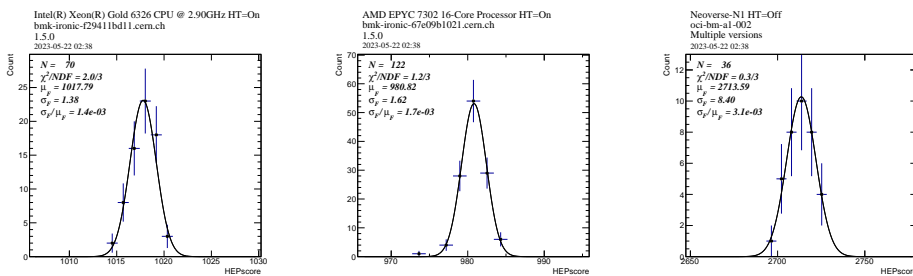


Figure 2. Histograms of the HEPsScore23 benchmark on an Intel, AMD and ARM processor. The fits to the histogram use a Gaussian distribution.

Each workload is encapsulated into a container with the software and input data needed to run the application. The software of the experiment is stored in the CVMFS file system [16] and then made accessible via a local folder inside the container. The set of containers of workloads is stored in a Gitlab repository at CERN. Each container includes a configuration file with a parameter for the number of events that allows one to adjust the duration of the execution (more details can be found in ref. [10]). Each workload is run three times and the geometric mean is taken as the benchmark (typically in units of events per second); this is identical to the method used for each workload component in HS06.

Each workload was validated on a set of dedicated servers at CERN to check the reliability and reproducibility, and in all cases, the results were found to be consistent at a level better than 1%. Once validated, the workloads were run on a large and diverse set of server systems provided by many WLCG sites. The benchmarks of the individual workloads from the server systems were compared to determine the correlations between applications. The time to run each workload and the studies of the correlations provided valuable input that was used to help select the workloads that are included in HEPsScore.

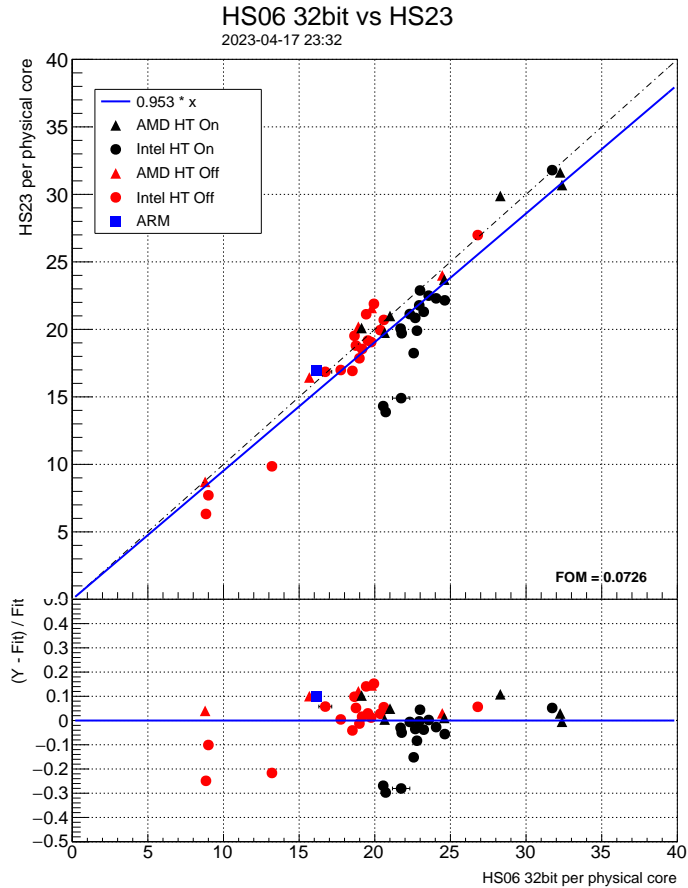


Figure 3. The blue line is a linear fit to the data points constrained to the origin and the dashed line has unity slope. The circles are measurements on Intel processors, the triangle points are measurements on AMD processors and the box point is an ARM processors measurement. The points in red (black) were taken with hyper-threading (on). The measurement on the single ARM processor is with hyper-threading off (ARM processors are not designed to operate with hyper-threading).

At the Benchmark Workshop, the criteria for selecting workload candidates for HEP-Score were discussed. The conclusion was that HEP-Score should be representative of the shares of computing usage of the experiments, it should run in a timely manner (3-6 hours) and provide complementary workloads (e.g. avoid the selection of highly correlated workloads). It is important that the workloads from the experiments be stable for a reasonable period (e.g. 2-3 years), and ideally, valid for at one LHC run period.

Seven workloads were selected to be part of the HEP-Score23 benchmark (the benchmark is called HEP-Score23 to indicate the year the benchmark was created). HEP-Score23 includes two workloads each from CMS (reconstruction and generation-simulation) and ATLAS (Sherpa-generation, reconstruction), and workloads from ALICE (digitization-reconstruction), Belle II (generation-simulation-reconstruction) and LHCb (simulation). The workloads were chosen to cover diverse types of applications, to be complementary in terms

of scaling across server types, to run in an acceptable time, and to represent complex event topologies. The time to run each workload ranges from 300 to 900 seconds on the reference server at CERN (Intel Xeon Gold 6326 CPU @ 2.90GHz). HEP-Score23 follows the methods used for HEP-SPEC06 of running each workload three times and taking the geometric mean of the three measurements. The total time to run the HEP-Score23 is approximately 3.5 hours on the reference server.

The workloads in HEP-Score23 are equally weighted. The Working Group studied different weighting schemes but found little difference. The choice of equal weighting of the workloads gives additional impact to ATLAS and CMS as they each provide two workloads.

As part of the validation process, the HEP-Score23 benchmark was measured on Intel, AMD and ARM servers at CERN. The results, shown in fig. 2, demonstrate the reproducibility of the HEP-Score23 benchmark on the different architectures to better than 1%. The HEP-Score23 benchmark was run on a wider set of servers with both hyper-threading on and off. In fig. 3, the HEP-Score23 benchmark is plotted against the HEP-SPEC06 (32-bit version).

In fig. 4, we show the ratio of the HEP-Score23 to HEP-SPEC06 as a function of the year in which the processor was introduced on the market (the colours of the points are identical to those used in fig. 3). The measurement of HEP-Score23 was normalized to the value HEP-SPEC06 on the reference machine (the data point for the reference machine is one of the points in 2021). It is observed that the ratio of HEP-Score23 to HEP-SPEC06 is less than unity for older machines and increases for newer servers.

4 Summary

The HEP-Score23 benchmark was released for use by WLCG sites in April 2023. HEP-Score23 is normalized to HEP-SPEC06 as measured on the reference machine to facilitate an easy transition for the sites and the WLCG accounting group. Sites are being asked to use HEP-Score23 (or both benchmarks) to evaluate newly procured hardware. Existing hardware will not need to be benchmarked with HEP-Score23 and sites can continue to use their current benchmarks based on HEP-SPEC06. The WLCG Management Board will review the situation and decide whether to make HEP-Score23 the required benchmark for future years.

References

- [1] WLCG, The Worldwide LHC Computing Grid. <http://wlcg.web.cern.ch>
- [2] HEPiX Benchmark Working group. <http://www.hepix.org/benchmarking.html>
- [3] Henning J.L., SPEC CPU2006 Benchmark Descriptions. SIGARCH Comput. Archit. News 34, 1. <https://doi.org/10.1145/1186736.1186737>
- [4] Michelotto M. et. al., A comparison of HEP code with SPEC benchmarks on multi-core worker nodes. J. Phys.: Conf. Ser. 219, 052009 (2010). <https://doi.org/10.1088/1742-6596/219/5/052009>.
- [5] Giordano D. and Santorinaiou E., Next Generation of HEP CPU benchmarks. J. Phys. Conf. Ser. 1525, 012073 (2020). <https://doi.org/10.1088/1742-6596/1525/1/012073>
- [6] SPEC CPU2017 Benchmark. <http://www.spec.org/cpu2017/>
- [7] Charpentier P., Benchmarking worker nodes using LHCb productions and comparing with HEP-SPEC06. J. Phys. Conf. Ser. 898, 082011 (2017). <https://doi.org/10.1088/1742-6596/898/8/082011>
- [8] Giordano D. et. al. HEPiX Benchmark Working Group status report. WLCG Workshop 2017, Manchester, UK. http://indico.cern.ch/event/609911/contributions/2620190/attachments/480455/2295576/WLCG_Workshop_2017_benchmarking_giordano.pdf

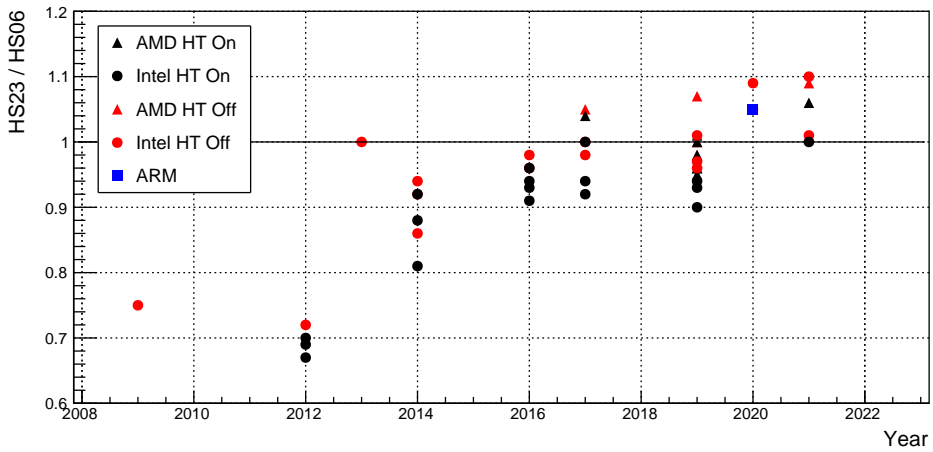


Figure 4. Plot of the ratio of the HEP-Score23 to HEP-SPEC06 benchmarks as a function of the year of the release of the server model. The circles are measurements on Intel processors, the triangle points are measurements on AMD processors and the box point are on ARM processors. The points in red (black) were taken with hyper-threading on). The measurement on the single ARM processor is with hyper-threading o (ARM processors are not designed to operate with hyper-threading).

[9] Valassi A. et. al. Using HEP experiment work ows for the benchmarking and accounting of WLCG computing resources. EPJ Web Conf 245:07035 (2020). <https://doi.org/10.1051/epjconf/202024507035>

[10] Giordano D.et. al., HEPiX Benchmarking Solution for WLCG Computing Resources. Comput Softw Big Sc5, 28 (2021). <https://doi.org/10.1007/s41781-021-00074-y>

[11] HEP-Score Workshop, September 2022, CERN. <https://indico.cern.ch/event/1170924/>

[12] Giordano D.et. al. The journey towards HEP-Score, the HEP-specific CPU benchmark for WLCG. To be published in the Proceedings of the 21st International Workshop on Advanced Computing and Analysis Techniques (ACAT) in Physics Research conference, Bari, Italy (2022). https://indico.cern.ch/event/1106990/contributions/4991202/attachments/25338174360309/ACAT_24_10_22_giordano.pdf.

[13] Sobie R.J.et. al. HEP-Score: a new CPU benchmark for the WLCG. WLCG Workshop 2022, Lancaster, UK. <https://indico.cern.ch/event/1162261/contributions/5092745/attachments/254384343802696/Sobie-WLCG-HEP-Score.pdf>

[14] Campana S. and Panzer B. A holistic study of the WLCG energy needs for HL-LHC. WLCG Workshop 2022, Lancaster, UK. <https://indico.cern.ch/event/1162261/contributions/5124364/attachments/25436874381025/WLCGEnergyNeeds.pdf>

[15] Simili E. et. al. Power Efficiency in HEP (a case between ARM and x86) To be published in the Proceedings of the 21st International Workshop on Advanced Computing and Analysis Techniques (ACAT) in Physics Research conference, Bari, Italy (2022). https://indico.cern.ch/event/1106990/contributions/4991256/attachments/253480143624687/PoW_ACAT2022.pdf

[16] Blomer J. et. al. New directions in the CernVM file system. J Phys Conf Ser 898:062031. <https://doi.org/10.1088/1742-6596/898/6/062031>